Figure 1: New version of Figure 2 in the original paper. As suggested by Reviewer fWzj, we replace the cube with tri-plane to represent the 3D face NeRF field. And change the "3D face" with "tri-plane" in the label to improve clarity.



(a) SegFormer-based Motion Adapter

(b) SegFormer Block in Motion Adapter

(c) Volume Renderer
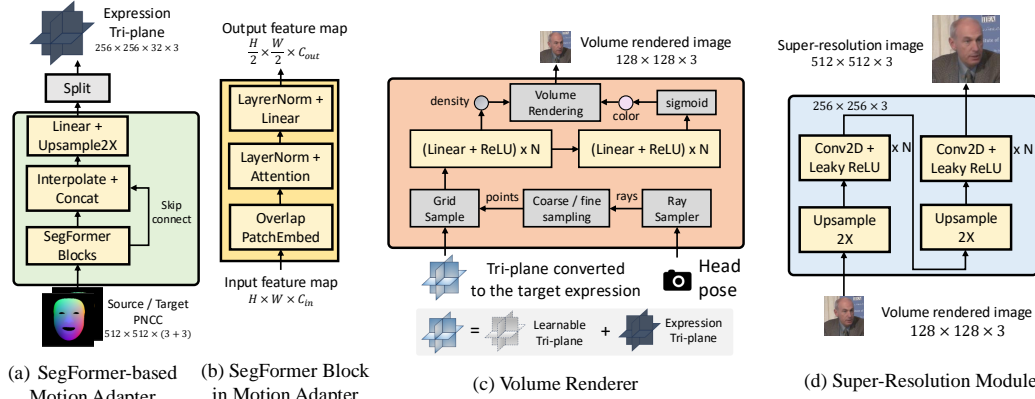
(d) Super-Resolution Module

Figure 2: Network details of each component in our generic model, as suggested by Reviewer fWzj.
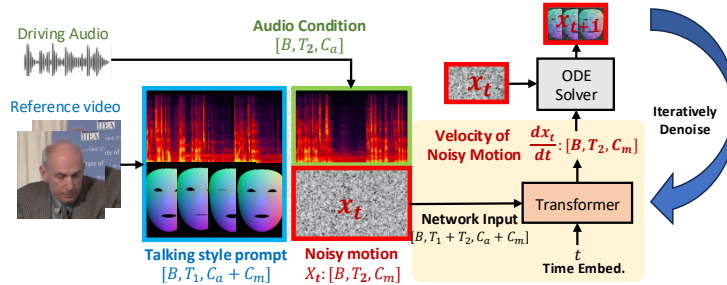


Figure 3: New version of Figure 3 in the original paper. As suggested by Reviewer fWzj, we improve the clarity of the inference process of ICS-A2M model.

Table 1: CMOS results on the lip-sync and expressiveness of various setting of our ICS-A2M model and the baseline ER-NeRF (ICCV 2023). CMOS score ranges from -3 to +3. Error bars are 95% confidence intervals.

| Settings | CMOS-lip-sync↑ | CMOS-expressive ↑ |
|---|---|---|
| #1. Ours (MimicTalk with ICS-A2M model) | **0.000** | **0.000** |
| #2. Ours w.o flow matching | $-0.438 \pm 0.273$ | $-0.895 \pm 0.387$ |
| #4. Ours w. style vector | $-0.361 \pm 0.204$ | $-0.532 \pm 0.210$ |
| #5. Ours w. style encoder | $-0.254 \pm 0.188$ | $-0.338 \pm 0.225$ |
| #6. Ours w.o sync loss | $-0.932 \pm 0.349$ | $-0.423 \pm 0.296$ |
| #6. ER-NeRF (ICCV 2023) | $-1.838 \pm 0.486$ | $-1.535 \pm 0.422$ |

Table 2: CMOS results on the style controllability and identity similarity of MimicTalk and StyleTalk. CMOS score ranges from -3 to +3. Error bars are 95% confidence intervals.

| Methods | CMOS-style-control↑ | CMOS-identitiy-similarity ↑ |
|---|---|---|
| #1. MimicTalk (ours) | **0.549 ± 0.225** | **1.735 ± 0.362** |
| #2. StyleTalk (AAAI 2023) | 0.000 | 0.000 |

1